

THE EQUATIONAL THEORIES PROJECT

MATTHEW BOLAN, JOSE BROX, MARIO CARNEIRO, MARTIN DVOŘÁK, ANDRÉS GOENS,
ZOLTAN KOCSIS, ALEX MEIBURG, PIETRO MONTICONE, JÉRÉMY SCANVIC, SHREYAS
SRINIVAS, TERENCE TAO, ANAND RAO TADIPATRI, VLAD TSARKLEVICH, DANIEL WEBER,
FAN ZHENG, ... (IN ALPHABETICAL ORDER)

ABSTRACT. We report on the *Equational Theories Project* (ETP), an online collaborative pilot project to explore new ways to collaborate in mathematics with machine assistance. The project sought to determine the implication graph between 4694 equational laws on magmas, by a combination of human-generated and automated proofs, all validated by the formal proof assistant language *Lean*. **state key outcomes**

CONTENTS

1. Introduction	3
1.1. Magmas and Equational Laws	3
1.2. The Equational Theories Project	4
1.3. Outcomes	5
1.4. Extensions	7
2. Notation and Mathematical Foundations	7
3. Formal Foundations	8
4. Project Management	9
4.1. Handling Scaling Issues	9
4.2. Other Design Considerations	10
4.3. Maintenance	10
5. Counterexample constructions	11
5.1. Finite magmas	11
5.2. Linear models	11

Date: December 4, 2024.

5.3.	Translation-invariant models	13
5.4.	The twisting semigroup	15
5.5.	Greedy constructions	16
5.6.	Modifying base models	18
5.7.	Ad hoc constructions	18
6.	Syntactic arguments	18
6.1.	Simple rewrites	18
6.2.	Matching invariants	19
6.3.	Confluence	20
6.4.	Complete rewriting systems	20
6.5.	Unique factorization	20
7.	Proof Automation	23
7.1.	Proof Techniques	24
7.2.	Proof Reconstruction and Integration	25
7.3.	Empirical Results	26
8.	Implications for Finite Magmas	27
9.	Order 5 laws	27
10.	Higman–Neumann laws	27
11.	AI and Machine Learning Contributions	27
11.1.	Graph ML: Directed link prediction on the implication graph	28
12.	User Interface	33
13.	Statistics and Experiments	33
14.	Data Management	33
15.	Reflections	33
16.	Conclusions and Future Directions	34

Acknowledgments	34
Appendix A. Numbering system	34
Appendix B. Author Contributions	37
References	38

1. INTRODUCTION

The purpose of this paper is to report on the *Equational Theories Project* (ETP)¹, a pilot project launched² in September 2024 to explore new ways to collaboratively work on mathematical research projects using machine assistance. The project goal, in the area of universal algebra, was selected to be particularly amenable to crowdsourced and computer-assisted techniques, while still being of mathematical research interest. **Describe outcomes**

1.1. Magmas and Equational Laws. In order to describe the mathematical goals of the ETP, we need some notation. A *magma* $M = (M, \diamond)$ is a set M (known as the *carrier*) together with a binary operation $\diamond: M \times M \rightarrow M$. An *equational law* for a magma, or *law* for short, is an identity involving \diamond and some formal indeterminates, which we will typically denote using the Roman letters x, y, z, w, u, v , as well as the formal equality symbol \simeq in place of the equality symbol $=$ to emphasize the formal nature of the law.

In the ETP, a unique number was assigned to each equational law, via a numbering system that we describe in Appendix A. For instance, the *commutative law* $x \diamond y \simeq y \diamond x$ is assigned the equation number Equation (E43), while the *associative law* $(x \diamond y) \diamond z \simeq x \diamond (y \diamond z)$ is assigned the equation number Equation (E4512). A list of all equations referred to by number in this paper is provided in Appendix A.

A magma M obeys a law E if the law E holds for all possible assignments of the indeterminate to M , in which case we write $M \models E$. Thus for instance $M \models E43$ if one has $x \diamond y = y \diamond x$ for all $x, y \in M$.

We say that a law E *entails* or *implies* another law E' if every magma that obeys E , also implies E' : $(M \models E) \implies (M \models E')$. We write this relation as $E \vdash E'$. We say that two laws are *equivalent* if they entail each other. For instance, the constant law $x \diamond y \simeq z \diamond w$ Equation (E46) can easily be seen to be equivalent to the law $x \diamond x \simeq y \diamond z$ Equation (E41). It is easy to see that \vdash is a pre-order, that is to say a partial order after one quotients by equivalence.

¹https://teorth.github.io/equational_theories/

²<https://terrytao.wordpress.com/2024/09/25>

In this entailment pre-ordering, the maximal element is given by the trivial law $x \simeq x$ Equation (E1), and the minimal element is given by the singleton law $x \simeq y$ Equation (E2), thus $E2 \vdash E \vdash E1$ for all laws E .

The *order* of an equational law is the number of occurrences of the magma operation. For instance, the commutative law Equation (E43) has order 2, while the associative law Equation (E4512) has order 4. We note some selected laws of small order that have previously appeared in the literature:

- The *central groupoid law* $x \simeq (y \diamond x) \diamond (x \diamond z)$ Equation (E168) is an order 3 law introduced by Evans [9] and studied further by Knuth [16] and many further authors, being closely related to central digraphs (also known as unique path property digraphs), and leading in particular to the discovery of the Knuth-Bendix algorithm [17]; see [22] for a more recent survey.
- *Tarski's axiom* $x \simeq y \diamond ((z \diamond (x \diamond (y \diamond z))))$ Equation (E543) is an order 4 law that was shown by Tarski [35] to characterize the operation of subtraction in an abelian group; that is to say, a magma M obeys Equation (E543) if and only if there is an abelian group structure on M for which $x \diamond y = x - y$ for all $x, y \in M$.
- In a similar vein, it was shown in [30] (see also [31]) that the order 4 law $x \simeq (y \diamond z) \diamond (y \diamond (x \diamond z))$ Equation (E1571) characterizes addition (or subtraction) in an abelian group of exponent 2; it was shown in [28] that the order 4 law $x \simeq (y \diamond ((x \diamond y) \diamond y)) \diamond (x \diamond (z \diamond y))$ Equation (E345169) characterizes the Sheffer stroke in a boolean algebra, and it was shown in [12] that the order 8 law $x \simeq y \diamond (((y \diamond y) \diamond x) \diamond z) \diamond (((y \diamond y) \diamond y) \diamond z)$ Equation (E42323216) characterizes division in a (not necessarily abelian) group.

Some further examples of laws characterizing well-known algebraic structures are listed in [27].

The Birkhoff completeness theorem [2, Th. 3.5.14] implies that an implication $E \vdash E'$ of equational laws holds if and only if the left-hand side of E' can be transformed into the right-hand side by a finite number of substitution rewrites using the law E . However, the problem of determining whether such an implication holds is undecidable in general [29]. Even when the order is small, some implications³ can require lengthy computer-assisted proofs; for instance, it was noted in [14] that the order 4 law $x \simeq (y \diamond x) \diamond ((x \diamond z) \diamond z)$ Equation (E1689) was equivalent to the singleton law Equation (E2), but all known proofs are computer-assisted.

1.2. The Equational Theories Project. As noted in Appendix A, there are 4694 equational laws of order at most 4. The primary mathematical goal of the ETP was to completely determine the *implication graph* for these laws, in which there is a directed edge from E to E' precisely when $E \vdash E'$. Such systematic determinations of implication graphs have been seen previously in the literature; for instance, in [32], the relations between 60 identities of Bol–Moufang type were established, and in the blog post [37, §17], some initial steps towards generating this graph for the first hundred or so laws on our list were performed. However,

³Another contemporaneous example of this phenomenon was the solution of the Robbins problem [26].

to our knowledge, the ETP is the first project to study such implications at the scale of thousands of laws.

The ETP requires the determination of the truth or falsity of $4694^2 = 22033636$ implications; while one can use properties such as the transitivity of entailment to reduce the work somewhat, this is clearly a task that requires significant automation. It was also a project highly amenable to crowdsourcing, in which different participants could work on developing different techniques, each of which could be used to fill out a different part of the implication graph. In this respect, the project could be compared with a Polymath project [11], which used online forums such as blogs and wikis to openly collaborate on a mathematical research problem. However, the Polymath model required human moderators to review and integrate the contributions of the participants, which clearly would not scale to the ETP which required the verification of over twenty million mathematical statements. Instead, the ETP was centered around a Github repository in which the formal mathematical contributions had to be entered in the proof assistant language *Lean*, where they could be automatically verified. In this respect, the ETP was more similar to the recently concluded Busy Beaver Challenge⁴, which was a similarly crowdsourced project that computed the fifth Busy Beaver number $BB(5)$ to be 47176870 through an analysis of 88664064 Turing machines, with the halting analysis being verified in a variety of computer languages, with the final formal proof written in the proof assistant language *Coq*. One of the aims of the ETP was to explore potential workflows for such collaborative, formally verified mathematical research projects that could serve as a model for future projects of this nature.

Secondary aims of the ETP included the possibility of discovering unusually interesting equational laws, or new experimental observations about such laws, that had not previously been noticed in the literature; and to develop benchmarks to assess the performance of automated theorem provers and other AI tools.

1.3. Outcomes. The ETP achieved its primary objective, with all of the implications formalized in the proof assistant language *Lean*, and can be found on the ETP GitHub repository. See Figure 1 for a small fragment of the implication graph produced. The experience of running such a large collaborative research project introduced several challenges, which we report upon in Section 4. Also, a variety of methods with varying degrees of automation or computer-assistance had to be developed to resolve all the implications, which had quite a variety of difficulty levels.

Of the 22033636 possible implications $E \vdash E'$, 8178279 (or 37.12%) would end up being true. To establish such positive implications $E \vdash E'$, the main techniques used were as follows:

- A very small number of positive implications were established and formalized by hand, mostly through direct rewriting of the laws; but this approach would not scale to the full project.
- Simple rewriting rules, for instance based on the observation that any law of the form $x \simeq f(y, z, \dots)$ was necessarily equivalent to the trivial law Equation (E2),

⁴<https://bbchallenge.org/>

- Greedy methods, in which either the multiplication table $(x, y) \mapsto x \diamond y$ or the function f determining a translation-invariant model are iteratively constructed by a greedy algorithm subject to a well-chosen ruleset, were effective in resolving many implications not easily disposed of by preceding methods. See Section 5.5.
- Starting with a simple base magma M obeying both E and E' , and either enlarging it to a larger magma $M' \supset M$, extending it to a magma N with a projection homomorphism $\pi : N \rightarrow M$, or modifying the multiplication table on a small number of values, also proved effective when combined with greedy methods. See Section 5.6.
- To each equation E one can associate a “twisting semigroup” S_E . If S_E is larger than $S_{E'}$, then this can often be used to disprove the implication $E \vdash E'$; see Section 5.4.
- Some *ad hoc* models based on existing mathematical objects, such as infinite trees, rings of polynomials, or “Kisielewicz models” utilizing the prime factorization of the natural numbers, could also handle some otherwise difficult cases. In some cases, the magma law induced some relevant and familiar structures, such as a directed graph or a partial order, which also helped guide counterexample constructions. See Section 5.7.
- Automated theorem provers were helpful in identifying which simplifying axioms could be added to the magma without jeopardizing the ability to refute the desired implication $E \vdash E'$.

1.4. **Extensions.** While the primary objective of the ETP was being completed, some additional related results were generated as spinoffs. Specifically:

- In Section 5.1 we report on a variant of the implication graph of the original set of 4692 equations, in which the magma is required to be finite.
- In Section 9 we report on classifying which of the 57882 distinct laws of order 5 are equivalent to the singleton law Equation (E2), either with or without the requirement that the magma be finite.
- In Section 10 we report on classifying the laws of order 8 that are equivalent to the Higman-Neumann law Equation (E42323216).

Also mention ML stuff, GUI

2. NOTATION AND MATHEMATICAL FOUNDATIONS

If M is a magma, we define the left and right multiplication operators $L_a, R_a : M \rightarrow M$ for $a \in M$ by the formula

$$(1) \quad L_x y = R_y x := x \diamond y.$$

We also define the squaring operator $S : M \rightarrow M$ by

$$(2) \quad Sx := x \diamond x = L_x x = R_x x$$

and the (right) cubing operator $C : M \rightarrow M$ by

$$(3) \quad Cx := Sx \diamond x = R_x^2 x.$$

If X is an alphabet, we let M_X denote the free magma generated by X , thus an element of M_X is either a letter in X , or of the form $w_1 \diamond w_2$ with $w_1, w_2 \in M_X$. Every function $f: X \rightarrow M$ into a magma M extends to a unique homomorphism $\varphi_f: M_X \rightarrow M$, that is to say a function for which $\varphi_f(w \diamond w') = \varphi_f(w) \diamond \varphi_f(w')$ for all $w, w' \in M_X$. Formally, an equational law with some indeterminates in X can be written as $w_1 \simeq w_2$ for some $w_1, w_2 \in M_X$; a magma M then obeys this law if and only if $\varphi_f(w_1) = \varphi_f(w_2)$ for all $f: X \rightarrow M$. We also define the order of a word $w \in M_X$ to be the number of occurrences of \diamond in the word, thus letters in X are of order 0, and the order of $w_1 \diamond w_2$ is the sum of the orders of w_1, w_2 .

A *theory* is a collection Γ of equational laws; we say that a magma M *satisfies* a theory, and write $M \models \Gamma$, if every law in Γ is obeyed by M . If E is an equational law, we write $\Gamma \vdash E$ if every magma that satisfies Γ also satisfies E . A *free magma* $M_{X,\Gamma}$ for such a theory Γ and an alphabet X is a magma satisfying Γ together with a map $\iota_{X,\Gamma}: X \rightarrow M_{X,\Gamma}$ which is universal in the sense that every function $f: X \rightarrow M$ to a magma M satisfying Γ uniquely determines a homomorphism $\varphi_{f,\Gamma}: M_{X,\Gamma} \rightarrow M$ such that $\varphi_{f,\Gamma} \circ \iota_{X,\Gamma} = f$. This magma is unique up to isomorphism; a canonical way to construct it is as the quotient M_X / \sim_Γ of the free magma M_X by the equivalence relation \sim_Γ given by declaring $w \sim_\Gamma w'$ if $\Gamma \vdash w \simeq w'$. If $\Gamma = \{E\}$ consists of a single law E , we write $M_{X,E}$, \sim_E , $\varphi_{f,E}$ for $M_{X,\{E\}}$, $\sim_{\{E\}}$, $\varphi_{f,\{E\}}$ respectively. **Give reference for free magmas relative to theories**

In general, the free magma $M_{X,\Gamma}$ is difficult to describe in a tractable form, but for some theories, one has a simple description:

Example 2.1 (Commutative and associative free magma). The free magma $M_{X,\{E43,E4512\}}$ for the commutative law Equation (E43) and the associative law Equation (E4512) is the free abelian semigroup generated by X (with $\iota_{X,\{E43,E4512\}}$ the obvious embedding map).

Example 2.2 (Left-absorptive free magma). The free magma $M_{X,\{E4\}}$ for the left-absorptive law Equation (E4) is the magma with carrier X and operation $x \diamond y = x$ (with $\iota_{X,E4}$ the identity).

Every magma M has an opposite M^{op} , which has the same carrier but the opposite operation $x \diamond^{\text{op}} y := y \diamond x$. A magma M obeys an equational law E if and only if its opposite M^{op} obeys the dual law E^* , defined by reversing the all operations. For instance, the dual of $x \diamond y \simeq x \diamond (y \diamond z)$ Equation (E327) is $y \diamond x \simeq (z \diamond y) \diamond x$, which in our numbering system we rewrite in normal form as $x \diamond y \simeq (z \diamond x) \diamond y$ Equation (E395).

We then see that the implication graph has a duality symmetry: given two equational laws E_1, E_2 , we have $E_1 \vdash E_2$ if and only if $E_1^* \vdash E_2^*$.

3. FORMAL FOUNDATIONS

TODO: expand this sketch.

Here we describe the Lean framework used to formalize the project, covering technical issues such as:

- Magma operation symbol issues
- Syntax ('LawX') versus semantics ('EquationX')
- "Universe hell" issues
- Additional verification (axiom checking, Leanchecker, etc.)
- Use of the 'conjecture' keyword
- Use of namespaces to avoid collisions between contributions. (Note: we messed up with this with FreeMagma! Should have namespaced back end results as well as front end ones.)
- Use of Facts command to efficiently handle large numbers of anti-implications at once

Upstream contributions:

- Mathlib contributions
- LeanBlueprint contributions

4. PROJECT MANAGEMENT

TODO: expand this sketch.

Shreyas Srinivas and Pietro Monticone have volunteered to take the lead on this section.

Discuss topics such as:

- Project generation from template
- Github issue management with labels and task management dashboard
- Continuous integration (builds, blueprint compilation, task status transition)
- Pre-push git hooks
- Use of blueprint (small note, see #406: blueprint chapters should be given names for stable URLs)
- Use of Lean Zulip (e.g. use of polls)

Maybe give some usage statistics, e.g. drawing from https://github.com/teorth/equational_theories/actions/metrics/usage

Mention that FLT is also using a similar workflow.

4.1. Handling Scaling Issues. Also mention some early human-managed efforts ("outstanding tasks", manually generated Hasse diagram, etc.) which suffices for the first one or two days of the project but rapidly became unable to handle the scale of the project.

Mention that some forethought in setting up a Github organizational structure with explicit admin roles etc. may have had some advantages if we had done so in the planning stages of the project, but it was workable without this structure (the main issue is that a single person - Terry - had to be the one to perform various technical admin actions).

Use of transitive reduction etc. to keep the Lean codebase manageable. Note that the project is large enough that one cannot simply accept arbitrary amounts of Lean code into the codebase, as this could make compilation times explode. Also note somewhere that transitive completion can be viewed as directed graph completion on a doubled graph consisting of laws and their formal negations.

Technical debt issues, e.g., complications stemming from an early decision to make `Equations.lean` and `AllEquations.lean` the ground truth of equations for other analysis and visualization tools, leading to the need to refactor when `AllEquations.lean` had to be split up for performance reasons.

Note that the "blueprint" that is now standard for guiding proof formalization projects is a bit too slow to keep up with this sort of project that is oriented instead about proving new results. Often new results are stated and formalized without passing through the blueprint, which is then either updated after the fact, or not at all. So the blueprint is more of an optional auxiliary guiding tool than an essential component of the workflow.

4.2. Other Design Considerations. Explain what "trusting" Lean really means in a large project. Highlight the kind of human issues that can interfere with this and how use of tools like external checkers and PR reviews by people maintaining the projects still matters. Provide guidelines on good practices (such as branch protection or watching out for spurious modifications in PRs, for example to the CI). Highlight the importance of following a proper process for discussing and accepting new tasks, avoiding overlaps etc. These issues are less likely to arise in projects with one clearly defined decision maker as in this case, and more likely to arise when the decision making has to be delegated to many maintainers.

Note that despite the guarantees provided by Lean, non-Lean components still contained bugs. For instance, an off-by-one error in an ATP run created a large number of spurious conjectures, and some early implementations of duality reductions (external to Lean) were similarly buggy. "Unit tests", e.g., checking conjectured outputs against Lean-validated outputs, or by theoretical results such as duality symmetry, were helpful, and the Equation Explorer visualization tool also helped human collaborators detect bugs.

Meta: documenting thoughts for the future record is quite valuable. It would be extremely tedious to try to reconstruct real-time impressions long after the fact just from the github commit history and Zulip chat archive.

4.3. Maintenance. Describe the role of maintainers and explain why they need to be conversant in the mathematics being formalised, as well as Lean. As such, the role of maintainers is often akin to postdocs or assistant profs in a research group who do some research of their own, but spend much of their time to guide others in performing their tasks, the key difference being that contributors are volunteers who choose their own tasks. Explain the tasks maintainers must perform. Examples:

- Reviewing proofs,
- Helping with proofs and theorem statements when people get stuck,

- Offering suggestions and guidance on how to produce shorter or more elegant proofs,
- Ensuring some basic standards are met in proof blocks that make proofs robust to upstream changes,
- Creating and maintaining CI processes,
- Responding to task requests,
- Evaluating theorem and definition formulations (for example unifying many theorem statements into one using FactsSyntax),
- Suggesting better ones where possible,
- Ensuring that there is no excessive and pointless overlap of content in different contributions (TODO: elaborate on what level of overlap was permissible and what we consider excessive).

5. COUNTEREXAMPLE CONSTRUCTIONS

In this section we collect the various techniques developed in the ETP to construct counterexamples to various implications $E \vdash E'$.

5.1. Finite magmas. **TODO: Expand this sketch**

Discuss semi-automated creation of finite counterexamples (as discussed here)

Describe various sources of example magmas used in counterexamples, including the ones listed here.

Also note some “negative results” - classes of finite magmas that did not yield many additional refutations, e.g. commutative 5x5 magmas.

Mention that many refutations are “immune” to finite models, point to Section 8 for details.

Using SAT solvers to find medium sized finite magmas obeying a given law? See this discussion.

Discuss computational and memory efficiencies needed to brute force over extremely large sets of magmas. SAT solving may be a better approach past a certain size!

5.2. Linear models. A fruitful source of counterexamples is the class of *linear magmas*, where the carrier M is a ring (which may be commutative or non-commutative, finite or infinite), and the operation \diamond is given by $x \diamond y = ax + by$ for some coefficients $a, b \in M$; one can also generalize this slightly to *affine magmas*, in which the operation is given by $x \diamond y = ax + by + c$, but for simplicity we shall focus on linear magmas here. It is easy to see that in a linear magma, any word $w(x_1, \dots, x_n)$ of n indeterminates also takes the linear form

$$w(x_1, \dots, x_n) = \sum_{i=1}^n P_{w,i}(a, b)x_i$$

for some (possibly non-commutative) polynomial $P_{w,i}$ in a, b with integer coefficients. Thus, a linear magma will obey an equational law $w_1 \simeq w_2$ if and only if the pair (a, b) lies in the (possibly non-commutative) variety

$$(4) \quad \{(a, b) \in M \times M : P_{w_1,i}(a, b) = P_{w_2,i}(a, b) \text{ for all } i\}.$$

As such, a necessary condition for such a law $w_1 \simeq w_2$ to entail another law $w'_1 \simeq w'_2$ is that one has the inclusion

$$\{(a, b) \in M \times M : P_{w_1,i}(a, b) = P_{w_2,i}(a, b) \text{ for all } i\} \subset \{(a, b) \in M \times M : P_{w'_1,i}(a, b) = P_{w'_2,i}(a, b) \text{ for all } i\}$$

for all rings M . For commutative rings, this criterion can be checked by standard Grobner basis techniques; in the noncommutative case one can use methods such as the diamond lemma [4].

Example 5.1 (Commutative counterexample). For the law $x = y \diamond (((x \diamond y) \diamond x) \diamond y)$ Equation (E1286), the variety Equation (4) can be computed to be

$$\{(a, b) \in M \times M : 1 = a + ba^3 + bab, 0 = a + ba^2b + b^2\}$$

while the variety for the idempotent law Equation (E3) is

$$\{(a, b) \in M : a + b = 1\}.$$

Thus to show that Equation (E1286) does not entail Equation (E3), it suffices to locate elements a, b of a ring M for one has $1 = a + ba^3 + bab$, $0 = a + ba^2b + b^2$, and $a + b \neq 1$. Here one can take a commutative example, for instance when $M = \mathbb{Z}/p\mathbb{Z}$ and $(p, a, b) = (11, 1, 7)$.

Example 5.2 (Noncommutative counterexample). For the law $x = y \diamond ((y \diamond (x \diamond z)) \diamond z)$ Equation (E1117), the variety Equation (4) can be computed to be

$$\{(a, b) \in M \times M : 1 = baba, 0 = a + ba^2, 0 = bab^2 + b^2\}$$

while the variety for $x = (x \diamond ((x \diamond x) \diamond x)) \diamond x$ Equation (E2441) is

$$\{(a, b) \in M \times M : a^2 + aba^2 + abab + ab^2 + b = 1\}.$$

Observe that if $ba = -1$, then (a, b) automatically lies in the first set, and lies in the second set if and only if $(ab + 1)(b - 1) = 0$. One can then show that Equation (E1117) does not imply Equation (E2441) by setting $a = L$, $b = -R$ where L, R are the left and right shift operators respectively on the ring of integer-valued sequences $\mathbb{Z}^{\mathbb{N}}$. With some *ad hoc* effort one can convert this example into a less linear, but simpler (and easier to formalize) example, namely the magma with carrier \mathbb{Z} and operation $x \diamond y = 2x - \lfloor y/2 \rfloor$.

Remark 5.3. As essentially observed in [1], if there is a commutative linear counterexample to an implication $E \vdash E'$, then by the Lefschetz principle this counterexample can be realized in a finite field \mathbb{F}_q for some prime power q (and by the Chebotarev density theorem one can in fact take q to be a prime, so that the carrier is of the form $\mathbb{Z}/p\mathbb{Z}$ for some prime p). As such, we have found that an effective way to refute implications by the commutative linear magma method is to simply perform a brute force search over linear magmas $x \diamond y = ax + by$ in $\mathbb{Z}/p\mathbb{Z}$ for various triples (p, a, b) . **Discuss performance of this method.**

On the other hand, the refutations obtained by non-commutative linear constructions need not have a finite model. For instance, consider the refutation $E1117 \not\models E2441$ from Example 5.2. The law Equation (E1117) can be rewritten as $L_y R_z L_y R_z x = x$. This implies that R_z is injective and L_y is surjective for all y, z . For finite magmas M , this then implies that

the L_y, R_z are in fact invertible, and hence we have also $R_z L_y R_z L_y x = x$, which implies Equation (E2441) by setting $x = y = z$. Thus the refutation $E1117 \not\vdash E2441$ is “immune” to finite counterexamples.

Remark 5.4. One can also consider nonlinear magma models, such as quadratic models $x \diamond y = ax^2 + bxy + cy^2 + dx + ey + f$ in a cyclic group $\mathbb{Z}/N\mathbb{Z}$. For small values of N , we have found such models somewhat useful in providing additional refutations of implications $E \vdash E'$ beyond what can be achieved by the linear or affine models. However, as the polynomials associated to a word $w(x_1, \dots, x_n)$ tend to be of high degree (exponential in the order of the word), it becomes quite rare for such models to obey a given equation E when N is large.

Remark 5.5. One can also consider the seemingly more general linear model $x \diamond y = ax + by$, where the carrier M is now an abelian group, and a, b act on M by homomorphisms, that is to say that they are elements of the endomorphism ring $\text{End}(M)$. However, this leads to exactly the same varieties Equation (4) (where M is now replaced by the endomorphism ring $\text{End}(M)$) and so does not increase the power of the linear model for the purposes of refuting implications.

Give some statistics of what proportion of refutations can be resolved by linear models.

On the other hand, there are certainly some refutations $E \not\vdash E'$ of implications that are “immune” to both commutative and non-commutative models, in the sense that all such models that obey E , also obey E' . One such example is the refutation $E1485 \vdash E151$, which we discuss further in Section 5.4 below.

5.3. Translation-invariant models. It is natural to look for counterexamples amongst magmas that obey a large number of symmetries. One such class of counterexamples are *translation-invariant models*, in which the carrier M is a group, and the left translations of this group form isomorphisms of the magma M . In the case of an abelian group $M = (M, +)$, such models take the form

$$(5) \quad x \diamond y = x + f(y - x)$$

for some function $f: M \rightarrow M$; in the case of a non-abelian group $M = (M, \cdot)$, such models instead take the form

$$(6) \quad x \diamond y = xf(x^{-1}y).$$

For such models, the verification of an equational law in n variables corresponds to a functional equation for f in $n - 1$ variables, as the translation symmetry allows one to normalize one variable to be the identity (say). This can simplify an implication to the point where an explicit counterexample can be found. These functional equations are trivial to analyze when $n = 1$. For $n = 2$, they are not as trivial, but still quite tractable, and has led to several refutations in practice. The method does not appear to be particularly effective for $n > 2$ due to the complexity of the functional equations.

Example 5.6 (Abelian example). For the law $x \simeq (x \diamond y) \diamond ((x \diamond y) \diamond y)$ Equation (E1648), we apply the abelian translation-invariant model Equation (5) with $y = x + h$ to obtain

$$\begin{aligned} x \diamond y &= x + f(h) \\ (x \diamond y) \diamond y &= x + f(h) + f(h - f(h)) \\ (x \diamond y) \diamond ((x \diamond y) \diamond y) &= x + f(h) + f(f(h - f(h))) \end{aligned}$$

so that this law obeys Equation (E1648) if and only if the functional equation

$$f(h) + f(f(h - f(h))) = 0$$

holds for all $h \in M$. Similarly, the law $x \simeq (x \diamond (x \diamond y)) \diamond y$ Equation (E206) is obeyed if and only if

$$f(f(h)) + f(h - f(f(h))) = 0$$

for all $h \in M$. One can now check that the function $f: \mathbb{Z} \rightarrow \mathbb{Z}$ defined by $f(h) := -\text{sgn}(h)$ (thus $f(h)$ equals -1 when h is positive, $+1$ when h is negative, and 0 when h is zero) obeys the first functional equation but not the second, thus establishing that $E1648 \not\models E206$.

Example 5.7 (Non-abelian example). We now obtain the opposite refutation $E206 \not\models E1648$ to Example 5.6 using the non-abelian translation-invariant model. By similar calculations to before, we now seek to find a function $f: M \rightarrow M$ on a non-abelian group (M, \cdot) that obeys the functional equation

$$(7) \quad f(f(h))f(f(f(h))^{-1}h) = 1$$

for all $h \in M$, but fails to obey the functional equation

$$(8) \quad f(h)f(f(f(h)^{-1}h)) = 1$$

for at least one $h \in M$. Now take M to be the group generated by three generators a, b, c subject to the relations $a^2 = b^2 = c^2 = 1$, or equivalently the group of reduced words in a, b, c with no adjacent letters in the word equal. We define

$$f(1) = 1, f(a) = b, f(b) = c, f(c) = a$$

and then $f(aw) = a$ for any non-empty reduced word w not starting with a , and similarly for b and c . The equation Equation (7) can be checked directly for $h = 1, a, b, c$. If $h = aw$ with w non-empty, reduced, and not starting with a , then $f(f(h))^{-1} = f(f(h)) = b$ and $f(f(f(h))^{-1}h) = f(baw) = b$, giving Equation (7) in this case, and similarly for cyclic permutations. Meanwhile, Equation (8) can be checked to fail for $h = a$.

Remark 5.8. The construction in Example 5.7 also has the following more “geometric” interpretation. The carrier M can be viewed as the infinite 3-regular tree, in which every vertex imposes a cyclic ordering on its 3 neighbors (for instance, if we embed M as a planar graph, we can use the clockwise ordering). For $x, y \in M$, we then define $x \diamond y$ to equal x if $x = y$. If y is instead a neighbor of x , we define $x \diamond y$ to be the next neighbor of x in the cyclic ordering. Finally, if y is distance two or more from x , we define $x \diamond y$ to be the neighbor of x that is closest to y . One can then check that this model obeys Equation (7) but not Equation (8).

Some refutations $E \not\models E'$ are “immune” by translation-invariant models, in that any translation-invariant model that obeys E , also obeys E' . One obstruction is that for such models, the squaring map S is necessarily an invertible map, since $Sx = x + f(0)$ in the abelian case

and $Sx = xf(1)$ in the non-abelian case. On the other hand, adding the assumption of invertibility of squares can sometimes make the implication $E \vdash E'$. For instance, the commutative law $x \diamond (y \diamond y) \simeq (y \diamond y) \diamond x$ Equation (E4482) for a square and an arbitrary element will imply the full commutative law Equation (E43) for translation-invariant models due to the surjectivity of S , but does not imply it in general (as one can easily see by considering models where S is constant).

5.4. The twisting semigroup. Suppose one has a magma M obeying a law E , that also enjoys some endomorphisms $T, U: M \rightarrow M$. Then one can “twist” the operation \diamond by T, U to obtain a new magma operation

$$(9) \quad x \diamond' y := Tx \diamond Uy.$$

If one then tests whether this new operation \diamond' obeys the same law E as the original operation \diamond , one will find that this will be the case provided that T, U obey a certain set of relations. The semigroup generated by formal generators T, U with these relations will be called the *twisting semigroup* Twist_E of E . This can be best illustrated with some examples.

Example 5.9. We compute the twisting semigroup of $x \simeq (y \diamond x) \diamond (x \diamond (z \diamond y))$ Equation (E1485). We test this law on the operation Equation (9), thus we consider whether

$$x = (y \diamond' x) \diamond' (x \diamond' (z \diamond' y))$$

holds for all $x, y, z \in M$. Substituting in Equation (9) and using the homomorphism property repeatedly, this reduces to

$$x = (T^2y \diamond T Ux) \diamond (UTx \diamond (U^2Tz \diamond U^3y)).$$

If we impose the conditions $TU = UT$, $T^2 = U^3$, then this equation would follow from Equation (E1485) (with x, y, z replaced with $T Ux$, T^2y , U^2Tz respectively). Thus the twisting semigroup Twist_{E1485} of Equation (E1485) is generated by two generators T, U subject to the relations $TU = UT = 1$, $T^2 = U^3$. This is a cyclic group of order 5, since the relations can be rewritten as $T^5 = 1$, $U = T^{-1}$.

Now consider $x \simeq (x \diamond x) \diamond (x \diamond x)$ Equation (E151). Applying the same procedure, we arrive at

$$x = (T^2x \diamond T Ux) \diamond (UTx \diamond U^2x)$$

so the twisting group Twist_{E151} is generated by two generators T, U subject to the relations $TU = UT = T^2 = U^2 = 1$. This is a cyclic group of order 2, since the relations can be rewritten as $T^2 = 1$, $U = T$.

Suppose the twisting semigroup Twist_E is not a quotient of $\text{Twist}_{E'}$, in the sense that the relations that define $\text{Twist}_{E'}$ are not obeyed by the generators of Twist_E . Then one can often disprove the implication $E \vdash E'$ by attempting the following procedure.

- First, locate a non-trivial magma M obeying the law E . Then the Cartesian power M^{Twist_E} of tuples $(x_W)_{W \in \text{Twist}_E}$, with the pointwise magma operation, will also obey E .

- Furthermore, this Cartesian power admits two endomorphisms T, U defined by

$$T(x_W)_{W \in \text{Twist}_E} = (x_{WT})_{W \in \text{Twist}_E}; U(x_W)_{W \in \text{Twist}_E} = (x_{WU})_{W \in \text{Twist}_E},$$

which obey the relations defining Twist_E .

- We now twist the magma operation \diamond on M^{Twist_E} by T, U to obtain a new magma operation \diamond' defined by Equation (9), that will still obey law E .
- Because T, U will not obey the relations defining $\text{Twist}_{E'}$, it is highly likely that this twisted operation will not obey E' , thus refuting the implication $E \vdash E'$.

For instance, a non-trivial finite model for Equation (E1485) is given by the finite field \mathbb{F}_2 of two elements with the NAND operation $x \diamond y := 1 - xy$. If we twist \mathbb{F}_2^5 by the left shift $T(x_i)_{i=1}^5 = (x_{i+1})_{i=1}^5$ and right shift $U(x_i)_{i=1}^5 = (x_{i-1})_{i=1}^5$, where we extend the indices periodically modulo 5, then the resulting operation

$$(x_i)_{i=1}^5 \diamond' (y_i)_{i=1}^5 := (1 - x_{i+1}y_{i-1})_{i=1}^5$$

on \mathbb{F}_2^5 will still obey Equation (E1485), but will not obey Equation (E151), thus showing that $E1485 \not\vdash E151$. This particular implication does not seem to be easily establishable by any of the other methods discussed in this paper.

Report on how large the twisting semigroups are in practice, and how many implications can be refuted by this method.

5.5. Greedy constructions. We have found *greedy extension methods*, or *greedy methods* for short, are a powerful way to refute implications, especially when the carrier M is allowed to be infinite. A basic implementation of this method is as follows. To build a magma operation $\diamond: M \times M \rightarrow M$ that obeys one law E but not another E' , one can first consider *partial magma operations* $\diamond: \Omega \rightarrow M$, defined on some subset Ω of $M \times M$. Thus $x \diamond y$ is defined if and only if $(x, y) \in \Omega$. A magma operation is then simply a partial operation which is *total* in the sense that $\Omega = M \times M$. We say that a partial magma operation is *finitely supported* if Ω is finite.

In the language of first-order logic, a partial magma operation can also be viewed as a ternary relation $R(x, y, z)$ on M with the axiom that $R(x, y, z) \text{ and } R(x, y, z') \implies z = z'$ for all $x, y, z \in M$. The support Ω is then the set of (x, y) for which $R(x, y, z)$ holds for some (necessarily unique) z , which one can then take to be the definition of $z = x \diamond y$.

We say that one partial operation $\diamond': \Omega' \rightarrow M$ *extends* another $\diamond: \Omega \rightarrow M$ if Ω' contains Ω , and $x \diamond y = x \diamond' y$ whenever $x \diamond y$ (and hence $x \diamond' y$) are defined. Given a sequence $\diamond_n: \Omega_n \rightarrow M$ of partial operations, each of which is an extension of the previous, we can define the *direct limit* $\diamond_\infty: \bigcup_n \Omega_n \rightarrow M$ to be the partial operation defined by $x \diamond_\infty y := x \diamond_n y$ whenever $(x, y) \in \Omega_n$.

Abstractly, the greedy algorithm strategy can now be described as follows.

Theorem 5.10 (Abstract greedy algorithm). *Let E, E' be equational laws, and let Γ be a theory of first-order sentences regarding a partial magma operations $\diamond: \Omega \rightarrow M$ on a carrier M . Assume the following axioms:*

- (i) (*Seed*) There exists a finitely supported partial magma operation $\diamond_0: \Omega_0 \rightarrow M$ satisfying Γ that contradicts E' , in the sense that there is some assignment of variables in E' in M such that both sides of E' are defined using \diamond_0 , but not equal to each other.
- (ii) (*Soundness*) If $\diamond_n: \Omega_n \rightarrow M$ is a sequence of partial magma operations obeying Γ with each \diamond_{n+1} an extension of \diamond_n , and the direct limit \diamond_∞ is total, then this limit obeys E .
- (iii) (*Greedy extension*) If $\diamond: \Omega \rightarrow M$ is a finitely supported partial magma operation obeying Γ , and $a, b \in M$, then there exists a finitely supported extension $\diamond': \Omega' \rightarrow M'$ of \diamond to a possibly larger carrier M' such that $a \diamond' b$ is defined.

Then $E \not\vdash E'$.

Proof. We work on the countably infinite carrier \mathbb{N} . By embedding the finitely supported operation \diamond_0 from axiom (i) into \mathbb{N} , we can assume without loss of generality that \diamond_0 has carrier \mathbb{N} . By similar relabeling, we can assume in (iii) that $M' = M$ when $M = \mathbb{N}$, since any elements of $M' \setminus \mathbb{N}$ that appear in Ω' can simply be reassigned to natural numbers that did not previously appear in Ω . We well-order the pairs in $\mathbb{N} \times \mathbb{N}$ by (a_n, b_n) for $n = 1, 2, \dots$. Iterating (iii) starting from \diamond_0 , we can thus create a sequence of finitely supported magma operations $\diamond_0, \diamond_1, \dots$ on \mathbb{N} obeying Γ , with each \diamond_{n+1} an extension of \diamond_n , and $a_n \diamond_n b_n$ defined for all $n \geq 1$. Then the direct limit \diamond_∞ of these operations is total, and does not obey E' thanks to axiom (i). On the other hand, by axiom (ii) it obeys E , and the claim follows. \square

We refer to Γ as the *rule set* for the greedy extension method. To apply Theorem 5.10 to obtain a refutation $E \vdash E'$, we have found the following trial-and-error method to work well in practice:

1. Start with a minimal rule set Γ that has just enough axioms to imply the soundness property for the given hypothesis E .
2. Attempt to establish the greedy extension property for this rule set by setting $a \diamond' b$ equal to a new element $c \notin M$, and then defining additional values of \diamond' as necessary to recover the axioms of Γ' .
3. If this can be done in all cases, then locate a seed \diamond_0 refuting the given target E' , and STOP.
4. If there is an obstruction (often due to a “collision” in which a given operation $x \diamond' y$ is required to equal two different values), add one or more rules to Γ to avoid this obstruction, and return to Step 2.

This method is not guaranteed to halt in finite time, as there may be increasingly lengthy sets of rules one has to add to Γ to avoid collisions. However, in practice we have found many of the refutations that could not be resolved by simpler techniques to be amenable to this method (or variants thereof, as discussed below).

One can automate the above procedure by using ATPs (or SAT solvers) to locate new rules that are necessary and sufficient resolve any potential collision (and which, *a posteriori*, can be seen to be necessarily consequences of the law E). **Describe performance of this automated method. Discuss the issue that some implications required a**

large SAT solver calculation that was difficult to formalize efficiently in Lean, prompting human-generated simplified proofs using smaller rulesets.

However, in some cases we have found it necessary to add “inspired” choices of rules that were not forced by the initial hypothesis E , but which simplified the analysis by removing problematic classes of collisions from consideration. We were unable to fully automate the process of guessing such choices; however, we found ATPs very useful for testing any proposed such guess. In particular, if an ATP was able to show that the existing ruleset, together with a proposed new rule A , implied E' , then this clearly indicated that one should not add A to the rule set Γ . Conversely, if an ATP failed to establish such an implication, this was evidence that this was a “safe” rule to impose.

We also found that human verification of the greedy extension property was a highly error-prone process, as the case analysis often included many delicate edge cases that were easy to overlook. Both ATPs and the Lean formalization therefore played a crucial role in verifying the human-written greedy arguments, often revealing important gaps in those arguments that required either minor or major revisions to the rule set.

Give examples

Discuss translation-invariant analogues

Discuss variants in which instead of using a single partial operation $\diamond : \Omega \rightarrow M$ or partial function f to capture the state of a greedy process, some more complicated structure is used instead. E.g. the Asterix construction.

Discuss the 1692 case in which one had to add an infinite tree of new values rather than a finite number.

5.6. Modifying base models. TODO: Expand this sketch

5.7. Ad hoc constructions. TODO: Expand this sketch

6. SYNTACTIC ARGUMENTS

Many proofs or refutations of implications (or equivalences) between two equational laws E, E' can be obtained from the syntactic form of the equation. We discuss some techniques here that were useful in the ETP.

6.1. Simple rewrites. Many equational laws E' can be formally deduced from a given law E by applying the Lean ‘rw’ tactic to rewrite E' repeatedly by some forward or backward application of E applied to arguments that match some portion of E . For instance, the commutative law Equation (E43) clearly implies $x \diamond (y \diamond z) \simeq (y \diamond z) \diamond x$ Equation (E4531) by a single such rewrite. A brute force application of such rewrite methods is already able to directly generate about 15,000 such implications, including many equivalences to the

singleton law Equation (E2) and the constant law Equation (E46). After applying transitive closure, this generates about four million further such implications.

A simple observation that already generates many equivalences is that any equation of the form $x \simeq f(y, z, \dots)$ necessarily is equivalent to the trivial law $x \simeq y$; similarly, an equation of the form $f(x, y) \simeq g(z, w, \dots)$ implies $f(x, y) \simeq f(x', y')$; and so forth. **Give some stats on how effective this is.**

6.2. Matching invariants. Fix an alphabet X . An *matching invariant* is an assignment $I: M_X \rightarrow \mathcal{I}$ of an object $I(w) \in \mathcal{I}$ in some space \mathcal{I} to each word $w \in M_X$ with the property that if an equational law $w_1 \simeq w_2$ has matching invariants $I(w_1) = I(w_2)$, then the same matching $I(w'_1) = I(w'_2)$ holds for any consequence $w'_1 \simeq w'_2$. In particular, if one law $I(w_1) = I(w_2)$ and $I(w'_1) \neq I(w'_2)$, then the law $w_1 \simeq w_2$ does not imply the law $w'_1 \simeq w'_2$.

A simple example of a matching invariant is the multiplicity $(n_x)_{x \in X}$ of variables of a word: if w_1, w_2 have all variables x appear the same number of times n_x in both words, then any rewriting of a word w using the law $w_1 \simeq w_2$ will preserve this property. Hence, if w'_1, w'_2 do not have the each variable appear the same number of times in both words, then $w_1 \simeq w_2$ cannot imply $w'_1 \simeq w'_2$. For instance, the commutative law Equation (E43) cannot imply the left-absorptive law Equation (E4).

One source of matching invariants comes from the free magma M_X of a theory:

Proposition 6.1 (Free magmas and matching invariants). *Let Γ be a theory, and let $\iota_{X,\Gamma}: X \rightarrow M_{X,\Gamma}$ be the map associated to the free magma $M_{X,\Gamma}$ for that theory. Then the map $I: M_X \rightarrow M_{X,\Gamma}$ defined by $I(w) := \varphi_{\iota_{X,\Gamma}}(w)$ is an invariant.*

Proof. Suppose that $w_1 \simeq w_2$ entails $w'_1 \simeq w'_2$, and that $I(w_1) = I(w_2)$. For any $f: X \rightarrow M_{X,\Gamma}$, the two maps $\varphi_f, \varphi_{f,\Gamma} \circ \varphi_{\iota_{X,\Gamma}}: M_X \rightarrow M_{X,\Gamma}$ are both homomorphisms that extend f , hence agree by the universal property of M_X , as displayed by the following commutative diagram:

$$\begin{array}{ccccc}
 & & X & & \\
 & \swarrow & \downarrow \iota_{X,\Gamma} & \searrow f & \\
 M_X & \xrightarrow{I = \varphi_{\iota_{X,\Gamma}}} & M_{X,\Gamma} & \xrightarrow{\varphi_{f,\Gamma}} & M_{X,\Gamma} \\
 & \searrow \varphi_f & & &
 \end{array}$$

In particular, the hypothesis $I(w_1) = I(w_2)$ implies that $\varphi_f(w_1) = \varphi_f(w_2)$ for all $f: X \rightarrow M_{X,\Gamma}$; that is to say, the magma $M_{X,\Gamma}$ obeys the law $w_1 \simeq w_2$, and hence also $w'_1 \simeq w'_2$ by hypothesis. In particular, $\varphi_{\iota_{X,\Gamma}}(w'_1) = \varphi_{\iota_{X,\Gamma}}(w'_2)$, which gives $I(w'_1) = I(w'_2)$ as required. \square

Example 6.2. If we take $\Gamma = \{E4\}$ to be the theory of the left-absorptive law Equation (E4) as described in Example 2.2, then the matching invariant $I(w)$ produced by Proposition 6.1 is the left-most letter of the alphabet X appearing in the word; for instance $I((x \diamond y) \diamond z) = x$. Thus, for example, the left-absorptive law Equation (E4) cannot imply the right-absorptive law Equation (E5).

Example 6.3. If we take $\Gamma = \{E43, E4512\}$ to be the theory of the commutative law Equation (E43) and the associative law Equation (E4512), then by Example 2.1, the associated invariant $I(w) = \sum_{x \in X} n_x e_x$ is the formal sum of all the generators e_x appearing in the word w , in the free abelian semigroup generated by those generators. This recovers the preceding observation that the multiplicities $(n_x)_{x \in X}$ form a matching invariant.

Example 6.4. Let $n \geq 1$ be a positive integer, and consider the theory $\Gamma = \{E43, E4512, E_n\}$ consisting of the previous theory $\{E43, E4512\}$ together with the order n law $L_x^y x = y$. One can check that the free magma $M_{X,\Gamma}$ can be described as the free group of exponent n with generators $e_x, x \in X$, with associated map $\iota_{X,\Gamma}: x \mapsto e_x$. The associated matching invariant $I(w) = \sum_{x \in X} n_x e_x$ is essentially the multiplicities $(n_x \bmod n)_{x \in X}$ modulo n , which gives a slightly stronger criterion than the preceding matching invariant for refuting implications. For example, the cubic idempotent law $x \simeq (x \diamond x) \diamond x$ Equation (E23) has matching invariants $e_x = 3e_x$ in the $n = 2$ case, and hence does not imply the idempotent law $x \simeq x \diamond x$ Equation (E3) since $e_x \neq 2e_x$ in the $n = 2$ case.

Give some statistics on how many refutations can be established by these methods.

Remark 6.5. One can also obtain matching invariants from the free objects associated to theories that involve additional operations beyond the magma operation \diamond , such as an identity element or an inverse operation. We leave the precise generalization of Proposition 6.1 to such theories to the interested reader.

6.3. Confluence. Define a confluent law and give some examples.

6.4. Complete rewriting systems. Define a complete rewriting system and give some examples.

6.5. Unique factorization. In general, the free magma $M_{X,E}$ for a given equational law E , which we can canonically define as M_X / \sim_E , is hard to describe explicitly; indeed, from the undecidability of implications, such a magma cannot be computably described for arbitrary E . Nevertheless, for some laws it is possible to obtain some partial understanding of $M_{X,E}$ from a syntactic perspective. For instance, if we can refute the equivalence $w'_1 \sim_E w'_2$ by constructing a counterexample magma M that obeys E but not $w'_1 \simeq w'_2$, then this implies that the representatives $\iota_{X,E}(w'_1), \iota_{X,E}(w'_2)$ of w'_1, w'_2 in $M_{X,E}$ are distinct.

We illustrate this approach with equations E of a left-absorptive form

$$(10) \quad x \simeq x \diamond f(x, y, z)$$

for some word $f(x, y, z)$, which imply the right-idempotent law Equation (E378).

An illustrative example is the law Equation (E854) depicted in Figure 1. Other examples are listed in Figure 2.

Lemma 6.6. *Equation (E854) is of the form Equation (10) and implies Equation (E378).*

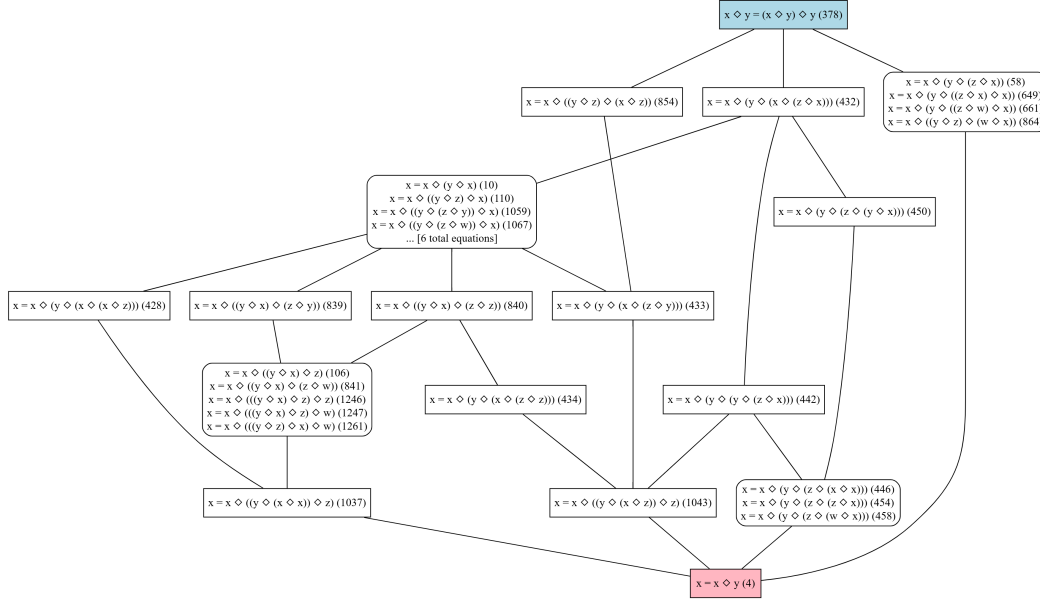


FIGURE 2. Equations similar to Equation (E854) that are of the form Equation (10) (possibly involving a fourth indeterminate w) and imply Equation (E378). For brevity, 70 equations equivalent to Equation (E4) have been omitted.

Proof. Clearly we have Equation (10) with $f(x, y, z) := (y \diamond z) \diamond (x \diamond z)$. From Equation (10) we have in any magma obeying Equation (E854) that

$$x = x \diamond f(x, S^2x, x) = x \diamond S(x \diamond S^2x) = x \diamond S(x \diamond f(x, x, x)) = x \diamond Sx.$$

This implies from a further application of Equation (10) that

$$y = y \diamond f(y, x, y) = (y \diamond Sy) \diamond ((x \diamond y) \diamond Sy) = f(x \diamond y, y, Sy)$$

and hence by Equation (10) again

$$(x \diamond y) \diamond y = x \diamond y$$

giving Equation (E378).

Let E be a law of the form Equation (10) that implies Equation (E378). We define a directed graph \rightarrow_E on words in M_X by declaring $w' \rightarrow_E w$ if $w \sim_E w'' \diamond w'$ for some $w' \in M_X$. By Equation (E378) (applied to the quotient magma $M_{X,E} = M_X / \sim_E$), this is equivalent to requiring that $w \sim_E w \diamond w'$. In particular, from Equation (10) we have $f(x, y, z) \rightarrow x$ for all x, y, z . Furthermore, the relation \rightarrow_E factors through \sim_E : if $w \sim_E \tilde{w}$ and $w' \sim_E \tilde{w}'$, then $w' \rightarrow_E w$ if and only if $\tilde{w}' \rightarrow_E \tilde{w}$.

Call a word $w \in M_X$ *irreducible* if it is not of the form $w = w_1 \diamond w_2$ with $w_2 \rightarrow_E w_1$. We can partially understand the equivalence relation \sim_E on irreducible words:

Theorem 6.7 (Description of equivalence). *Let E be an equation of the form Equation (10). Let w be an irreducible word, and let w' be a word with $w \sim_E w'$.*

(i) If w is a product $w = w_1 \diamond w_2$, then w' takes the form

$$w' = (((w'_1 \diamond w'_2) \diamond v_1) \diamond \dots \diamond v_n)$$

for some $w'_1 \sim_E w_1$, $w'_2 \sim_E w_2$, some $n \geq 0$, and some words v_1, \dots, v_n such that for all $0 \leq i < n$, v_{i+1} is of the form

$$v_{i+1} \sim_E f(x_i, y_i, z_i)$$

for some x_i, y_i, z_i with

$$x_i \sim_E (((w'_1 \diamond w'_2) \diamond v_1) \diamond \dots \diamond v_i).$$

In particular, $v_{i+1} \rightarrow_E x_i$.

(ii) Similarly, if $w \in X$ is a generator of M_X , then w' takes the form

$$w' = ((w \diamond v_1) \diamond \dots \diamond v_n)$$

for some $n \geq 0$, and some words v_1, \dots, v_n such that for all $0 \leq i < n$, v_{i+1} is of the form

$$v_{i+1} \sim_E f(x_i, y_i, z_i)$$

for some x_i, y_i, z_i with

$$x_i \sim_E ((w \diamond v_1) \diamond \dots \diamond v_i).$$

In particular, $v_{i+1} \rightarrow_E x_i$.

Conversely, any word of the above forms is equivalent to w .

Proof. We just verify claim (i), as claim (ii) is similar. The converse direction is clear from Equation (10) (after quotienting by \sim_E), so it suffices to prove the forward claim. By the Birkhoff completeness theorem, it suffices to prove that the class of words described by (i) is preserved by any term rewriting operation, in which a term in the word is replaced by an equivalent term using Equation (10). Clearly the term being rewritten is in w'_1 or w'_2 then the form of the word is preserved, and similarly if the term being rewritten is in one of the v_i . The only remaining case is if we are rewriting a term of the form

$$x_i = (((w'_1 \diamond w'_2) \diamond v_1) \diamond \dots \diamond v_i).$$

If $i > 0$ we can rewrite this term down to x_{i-1} , and this still preserves the required form (decrementing n by one). If $i = 0$ then we cannot perform such a rewriting because of the irreducibility of $w_1 \diamond w_2$ and hence $w'_1 \diamond w'_2$. Finally, we can rewrite x_i to $x_i \diamond v$ where v is of the form

$$v_i = f(x_i, y, z),$$

and after some relabeling we are again of the required form (now incrementing n by one). This covers all possible term rewriting operations, giving the claim. \square

Specializing to the case where w, w' are both irreducible, we conclude

Corollary 6.8 (Unique factorization). *Two irreducible words w, w' are equivalent if and only if they are either the same generator of X , or are of the form $w = w_1 \diamond w_2$, $w' = w'_1 \diamond w'_2$ with $w_1 \sim_E w'_1$ and $w_2 \sim_E w'_2$.*

As an application of this corollary, we establish

Proposition 6.9 (E854 does not imply E3316). *Equation Equation (E854) does not imply Equation (E3316).*

Proof. (Sketch) We work in the free group M_X on two generators $X = \{x, y\}$. It suffices to show that

$$x \diamond y \not\sim_{E854} x \diamond (y \diamond (x \diamond y)).$$

Suppose this were not the case, then by Corollary 6.8 one of the following statements must hold:

- (i) $y \rightarrow_{E854} x$.
- (ii) $(y \diamond (x \diamond y)) \rightarrow_{E854} x$.
- (iii) $y \diamond (x \diamond y) \sim_{E854} y$.

If (i) holds, then we have $x \diamond y = x$ must hold in M_X / \sim_E , hence Equation (E854) would imply Equation (E4). However, it is possible to refute this implication by a finite counterexample.

Similarly, if (iii) held, then Equation (E854) would have to imply Equation (E10), but this can also be refuted in a finite magma.

Finally, if (ii) held, then the claim

$$x \diamond y \sim x \diamond (y \diamond (x \diamond y))$$

to refute simplifies to

$$x \diamond y \sim x$$

and we are back to (i), which we already know not to be the case. \square

7. PROOF AUTOMATION

In this project we used proof automation in two ways: automated theorem provers (ATPs) and Lean tactics. ATPs are generally stand-alone tools that implement a (semi-) decision procedure for a given formal language or related set of languages. For example, Vampire [20] is an ATP focused primarily on first-order logic using superposition, which we used extensively in this project. ATPs are generally complex pieces of software, and can often have bugs. While we could choose to trust the results of an ATP, and in this project we don't do so either. Instead, we use proof certificates, which many ATPs can produce, to reconstruct a proof in Lean. This process depends on the proof certificate and the ATP, and we describe it for the main reconstruction we've done.

Tactics in Lean, on the other hand, are meta-programs [8] that builds proofs. In other words, tactics are programs that operate at the meta-level of Lean code: they essentially take in Lean code as input and produce Lean code as output. In this manner, they look like another keyword in the language, and are tightly integrated by producing proofs directly. Under the hood, they can implement essentially anything, from syntactic sugar to full decision

procedures. The **duper** tactic [5], for example, implements a superposition calculus, similar Vampire’s, but for dependent types — Lean’s underlying logical foundation.

In the rest of this section we describe the different proof automation techniques and ATPs/tactics we used. We first discuss the different proof methods used: primarily superposition and equational reasoning, we then discuss the integration in Lean, and finally we report some basic empirical results from this project.

7.1. Proof Techniques. The main two families of ATPs and tactics we used here are superposition/saturation-based and equational reasoning ones. In this context we also include SMT solvers, which combine specific decision procedures for theories, like congruence closure for equational reasoning, with satisfiability (SAT) solving [?]. Finally, we also used **aesop** [24], which implements a version of tableau search. This was used mainly to help specific constructions in refutations, and is not specific to proving or disproving magma implications in this sense. We will describe our use of **aesop** more in Section 7.2 below.

Saturation. Most of the ATPs we used extensively in this project rely primarily on saturation procedures in the superposition calculus. This is the case for Vampire [20], and the paper also gives a more in-depth explanation of this. See also [3] for a gentler exposition. The core idea of these provers is that they take a set of assumptions and conjecture, expressed in — say — first-order logic. They take the conjecture and negate it, adding this negation to the set of assumptions, which are all put in some normal form. The ATP then try to refute the negation by applying rules of the underlying calculus, until they find a proof of false (a contradiction). In this case, the conjecture was (classically) true, and the ATP has found a proof by contradiction, often called a “refutation” or “saturation” proof.

The underlying calculus varies from system to system, but they often have a variant of a resolution clause, a clause of a form:

$$\frac{C \vee L \quad D \vee \neg L}{C \vee D}$$

This can also be read as $C \vee L \quad D \vee \neg L$ implies $C \vee D$, where C, D, L are formulas in e.g. first-order logic. Superposition calculi have a variant of this rule that deals with equality directly, and thus are more efficient at reasoning about equality.

In this project we used Vampire [20], Duper [5] and Prover9 and Mace4 [25] which are all based on variants of saturation for proving.

Equational Reasoning. Equational reasoning is a type of reasoning based on equational logic and rewriting with congruence [2], see Chapter ?? for a discussion of its foundations in universal algebra. In general, it takes a series of equations and determines whether another equation can be deduced from it. A core tool in equational reasoning are e-graphs, a data structure used to represent equivalence classes of terms. The procedure of congruence closure that can be used to decide ground equational problems (i.e. problems without variables) and as a semi-decision procedure in general.

SMT solvers like Z3 [7] use equational reasoning for deciding the theory of equality with uninterpreted functions [21, 6]. On the other hand, equality saturation [36] uses e-graphs by extending congruence closure to a more controlled search, enabling optimization and conditional rewriting. One of the main advantages of using equational reasoning to reason about implications of magma laws is that we get very explicit proofs: a proof that $l \vdash l'$ is given by a sequence of rewrites that starts at the left-hand side of l' and arrives at the right-hand side through applications of l .

In this project we used Z3 [7], Prover9 and Mace4 [25], a custom ATP for magmas based on egg [36], and the Lean egg tactic [18, 33], which all work with equational logic. We have also reasoned with manual (custom written) heuristics about simple rewrites.

7.2. Proof Reconstruction and Integration. While ATPs are very useful to solve questions for this project, they generally don't integrate well with Lean (except for Lean tactics like `duper`, `aesop` and `egg`). A bug in an ATP could lead to an unsound proof, or worse, an incorrect result. To avoid this, and having to trust large codebases from ATPs, we take results found by the different provers and (re-)construct proofs in Lean.

Integration implies two steps: invoking the decision procedures/ATPs (translating the problem from Lean into the languages and logics they use), and conversely, (re-)constructing the results from the decision procedure as a (persistent) Lean proof. These two aspects present different challenges, and require different strategies, depending mostly on the kind of proof strategy the decision procedure uses.

For saturation proofs ... (TODO: add explanations from <https://leanprover.zulipchat.com/#narrow/channel/458659-Equational/topic/Vampire.20-.3E.20Lean/near/478147138>)

For equational proofs from external provers (e.g. MagmaEgg), we used also a simple version of reconstruction by (re-)constructing the proofs of equality from an explanation, using congruence lemmas from Lean. In the case of equational proofs from the `egg` tactic, these could be converted into a series of `calc` steps, documenting the explicit calculation, using the `calcify` tactic⁵.

In general, we've observed that there are multiple ways of integrating decisions procedures within Lean, with different levels of integration.

- (1) Using a Lean tactic, which calls a decision procedure written in Lean (like `aesop` or `duper`).
- (2) Using a Lean tactic, which calls an existing (external) ATP and reconstructs a proof term from the ATP's result (like `bv_decide` or `egg`).
- (3) Using an external script which calls an existing ATP and generates a source file `.lean` which captures the result explicitly.

This project primarily used the least integrated approach, (Option 3), as it was fastest and required no dependencies on the other contributors. This also has drawbacks: primarily,

⁵<https://github.com/nomeata/lean-calcify/>

while the upfront effort is lower, the effort to use is higher than with a tactic, once the tactic is developed. It also makes integration with larger code bases more difficult: in this project the (mathematical) dependencies were by design very minimal, Magmas are simple, and we built our own definitions for them. For example, integration with the typeclass system become much more difficult when working with more complex mathematical objects that build upon many layers of structure.

Semi-Automated Counterexample Guidance. Another use of ATPs has been in a semi-automatic fashion, to find counter-examples. The general strategy was to use ATPs to find counter-examples to implications by building magmas iteratively. If we want to build a counterexample to $l \vdash l'$, we want to build construct a magma where l holds but l' doesn't. In this method, we iteratively strengthen a construction with additional hypotheses, and use the ATP to check whether these hypotheses are not too strong (to imply l') or unsound (to disallow l).

TODO: this should also be expanded more, at least with references to some of the constructions in other chapters.

While equational theories can also be used in a semi-automatic fashion to prove equations [18], the positive implications in this project were all simple enough that we did not need a semi-automatic approach for them.

7.3. Empirical Results. Finally, we report some empirical results from use of ATPs for this project, in terms of performance. The aim of this section is not to be a careful evaluation and benchmark comparison of the different ATPs; instead, we present our work here as a more informal “field report” documenting our experiences. In particular, we don't believe we can draw firm conclusions about the overall capabilities of the different ATPs. Rather, this serves as a use-case documenting the experience of (mostly) novice users.

TODO: throw a couple of “benchmarking” tables for the same ATP with different parameters and for different ATPs, talk about some relative gains in time (changing parameters we saw a 500 times speedup on this particular problem), etc. This is knowledge I think we have gained to some extent, and certainly I would have been glad to receive this kind of hints before we started!”. Then leave it as an interesting open problem to properly develop and measure benchmarks for ATPs based on this project.

TODO: Any comparative study of semi-automated methods with fully automated ones? In principle, the semi-automated approach could be more automated using a script or “agent” to call various theorem provers. See this discussion.

Note: when evaluating the performance of any particular automated implication tool, we should do a fresh run on the entire base of implications, rather than rely on whatever implications produced by the tool remain in the Lean codebase by the time of writing the paper. This is because (a) many of the previous runs focused only on those implications that were not already ruled out by earlier contributions, and (b) some pruning has been applied subsequent to the initial runs to improve compilation performance. Of course, these

runs would not need to be added to the (presumably optimized) codebase at that point, but would be purely for the purpose of gathering performance statistics. More discussion here.

See this discussion on the value of using different ATPs and setting run time parameters etc. at different values.

What are the hardest implications to prove? See this discussion.

Make a note of the possible alternate strategy to prove implications outlined here.

8. IMPLICATIONS FOR FINITE MAGMAS

TODO: Expand this sketch

Introduce Austin pairs.

Recap discussion from <https://leanprover.zulipchat.com/#narrow/channel/458659-Equational/topic/Austin.20pairs>

9. ORDER 5 LAWS

TODO: report on laws on order 5

10. HIGMAN–NEUMANN LAWS

TODO: report on Higman–Neumann laws

11. AI AND MACHINE LEARNING CONTRIBUTIONS

TODO: expand this sketch

- Claude assistance in coding front ends.
- ChatGPT to guess a complete rewriting system.
- Minor use of GitHub Copilot to autocomplete code in Lean and other languages.
- See this discussion.
- Directed Link Prediction (DLP) on the implication graph with multiple GNN autoencoders (see related Zulip topic).

ML experiments to learn the implication graph.

11.1. Graph ML: Directed link prediction on the implication graph. We experimented with various Graph Neural Network (GNN) autoencoders to predict missing edges, providing a way to estimate the truth value of unproven implications. To assess these models, we defined three test sets focusing on edge existence, directionality, and bidirectionality. The results give insight into how these models handle dense, directed graphs like ours. A more detailed report follows below.

11.1.1. Motivation. Directed Link Prediction [13] is a method enabling machine learning and deep learning models to predict missing edges in a directed graph. For our implication graph, this translates to predicting the truth values of unproven implications. This task serves as a necessary first step for advancing in the following directions:

- (1) **Reasoning over mathematical knowledge graphs:** Recent advancements allow language models to integrate information from multiple modalities. For example, [39, 38] share information between corresponding layers of Language Models (LMs) and Graph Neural Networks (GNNs), enabling simultaneous learning from text corpora and graph-structured data within the same expert domain. By leveraging both modalities, the language model can better *structure* its knowledge and respond to complex queries. This dual learning process combines masked language modeling for text with link prediction on the graph, highlighting the importance of link prediction for robust reasoning.
- (2) **Higher-order implication graphs:** Our implication graph currently represents only implications of the form $p \implies q$, not more complex ones like $(p \wedge r) \implies q$. Extending to such higher-order edges would likely involve connecting sets of nodes, thereby requiring hypergraph representations. For a systematic overview, see [15]. While specific hypergraph neural architectures exist [10], we believe it is still conceptually important to apply Directed Link Prediction to simpler implication graphs first, providing insights and guiding principles that can anticipate challenges in higher-order graph representation learning.

11.1.2. Data. We used this edge list, generated on October 20, 2024, with the following commands:

```
lake exe extract_implications outcomes > data/tmp/outcomes.json
scripts/generate_edgelist_csv.py
```

The structure of the implication graph is summarized below:

```
graph_summary = {
    "total_nodes": 4694,
    "total_directed_edges": 8178279,
    "edge_density_percentage": 37, # Percentage of possible edges that exist
    "bidirectional_edges": 2475610,
    "bidirectional_percentage": 30 # Percentage of all edges that are bidirectional
}
```

Below is a summary of the edge types in the graph:

```
edge_counts = {
    "explicit_conjecture_false": 92,
    "explicit_proof_false": 582316,
    "explicit_proof_true": 10657,
    "implicit_conjecture_false": 142,
    "implicit_proof_false": 13272681,
    "implicit_proof_true": 8167622,
    "unknown": 126
}
```

Edges are labeled according to the following scheme:

```
edge_labels = {
    "implicit_proof_true": 1,
    "explicit_proof_false": 0,
    "implicit_proof_false": 0,
    "explicit_proof_true": 1,
    "explicit_conjecture_false": 0,
    "implicit_conjecture_false": 0,
    "unknown": 0
}
```

The **unknown** class contains a very small number of edges (126), so their impact on the training phase is expected to be negligible. Future approaches might address this class by excluding **unknown** edges from the training set.

11.1.3. *Methods.* Consider a directed graph $G = (V, E)$ where $E = \{(u, v) \mid u, v \in V\}$ is the edge set and $|V| = n$. We assume that each node is associated with a feature vector, resulting in an $X \in \mathbb{R}^{n \times f}$ feature matrix.

We define the existing edges $(a, b) \in E$ as *positives* and the non-existing edges $(c, d) \notin E$ as *negatives*.

Intuitively, performing Directed Link Prediction (DLP) on G involves randomly splitting E into three disjoint sets: E_{train} , E_{val} , and E_{test} , such that:

- E_{train} is the training set,
- E_{val} is the validation set,
- E_{test} is the test set,
- and $E = E_{\text{train}} \dot{\cup} E_{\text{val}} \dot{\cup} E_{\text{test}}$.

The model then learns from $G_{\text{train}} = (V, E_{\text{train}})$ to map the topological and feature-related information of two nodes u and v to a probability p_{uv} that $(u, v) \in E_{\text{test}}$.

However, this setup presents two key issues among others:

- (1) The model learns only from *positives*, so it cannot recognize *negatives*.
- (2) The model is evaluated only on *positives*, preventing us from measuring its ability to identify *negatives*.

To address these limitations, we adopted the setup proposed in [34]. Specifically, we redefined E_{train} to E_{train}^p (*positives*) and introduced:

$$E_{\text{train}} = E_{\text{train}}^p \dot{\cup} E_{\text{train}}^n$$

where E_{train}^n includes all possible *negatives* in $G_{\text{train}} = (V, E_{\text{train}}^p)$. The model is now required to predict the non-existence of edges in E_{train}^n .

Similarly, if we redefine E_{test} as follows:

$$E_{\text{test}} = E_{\text{test}}^p \dot{\cup} E_{\text{test}}^n$$

where E_{test}^n is a *random* sample of *negatives*, the model's evaluation would fail to capture two crucial aspects:

- (1) The model's ability to distinguish (u, v) from (v, u) for all $(u, v) \in E_{\text{test}}^p$.
- (2) The model's ability to identify bi-implications.

These limitations arise from the random selection of negative edges in E_{test}^n . To address this, we define three distinct test sets: E_{test}^G , E_{test}^D , and E_{test}^B , to evaluate different facets of the model's performance:

- **General Test Set** (E_{test}^G): Here, $E_{\text{test}} = E_{\text{test}}^p \dot{\cup} E_{\text{test}}^n$, where E_{test}^n is a random sample of non-existent edges with the same cardinality as E_{test}^p . This set assesses the model's ability to detect the presence of edges, regardless of direction. A model that cannot distinguish edge direction may still perform well on this set, highlighting the need for the following two additional test sets.
- **Directional Test Set** (E_{test}^D): Defined as $E_{\text{test}}^{\text{up}} \dot{\cup} \tilde{E}_{\text{test}}^{\text{up}}$, where $E_{\text{test}}^{\text{up}}$ consists of unidirectional edges in E_{test}^p , and $\tilde{E}_{\text{test}}^{\text{up}}$ contains their reverses (negatives by construction). This set evaluates the model's ability to recognize edge direction, making it suitable for assessing direction-sensitive models.
- **Bidirectional Test Set** (E_{test}^B): Defined as $E_{\text{test}}^{\text{bp}} \dot{\cup} E_{\text{test}}^{\text{n}}$, where $E_{\text{test}}^{\text{bp}}$ contains one direction of all bidirectional edges in E_{test}^p , and $E_{\text{test}}^{\text{n}} \subset \tilde{E}$ includes a subset of their reverses. This set evaluates the model's ability to identify bi-implications, as each edge in E_{test}^B has a reverse that is positive, but only half are bidirectional in practice.

We tested the following models:

- **GAE** [13]
- **Gravity-GAE** [34]
- **Source/Target-GAE** [34]
- **DiGAE** [19]
- **MagNet** [40]

All these models are graph-based autoencoders with distinct encoder-decoder architectures. Notably, GAE is the only model structurally unable to differentiate edge directions. Each model outputs the probability that an ordered pair of nodes has a directed edge between them, with nodes represented using one-hot encodings as features.

We trained the models using Binary Cross Entropy as the loss function, balancing the contribution of positive and negative edges through re-weighting. On the *General* test set, we evaluated the following metrics:

- **AUC** (Area Under the ROC Curve): Measures the probability that the model ranks a random positive edge higher than a random negative edge. Higher values indicate better discrimination between positive and negative edges.
- **AUPRC** (Area Under Precision-Recall Curve): Assesses model performance, particularly in cases of class imbalance. AUPRC is more robust to imbalanced data than AUC.
- **Hits@K**: Evaluates the fraction of times a positive edge is ranked within the top K positions among personalized negative samples [23]. Briefly, given a positive edge, its M personalized negative samples are M negative edges with the same head but different tails. We calculate Hits@K for $K = 1, 3, 10$ to assess ranking quality, and set $M = 100$. Therefore, Hits@K = 0.1 means that on average, the model ranks a positive edge within the highest-ranked K personalized negatives 10% of the time.
- **MRR** (Mean Reciprocal Rank): Computes the average reciprocal rank of positive edges among their personalized negative samples [23] (the same as those used for Hits@K) providing an overall measure of ranking accuracy. For instance, $MRR = 0.1$ means that on average, the model ranks a positive edge as 10th among M personalized negative samples.

Each metric ranges from 0 to 1, with higher values reflecting improved performance. Based on prior work, we expect AUC and AUPRC scores to approach 1, while MRR and Hits@K often yield values around 0.15 for similar undirected tasks [23]. *Directional* and *Bidirectional* performances were evaluated using only AUC and AUPRC, since Hits@K and MRR are hardly definable in those scenarios. The number of training epochs was optimized through Early Stopping on the *General* validation set performance (given by the sum of AUC and AUPRC).

11.1.4. *Results.* The results below represent average performance over six random splits of E_{train} , E_{val} , and E_{test} while keeping the model’s seed fixed for fair comparison. The *training* / *validation* / *test* split proportions are:

- 85/5/10 for unidirectional edges,

- 65/15/30 for bidirectional edges.

Model	G_ROC_AUC	G_AUPRC	G_Hits@1	G_Hits@3
gae	$0.8484 \pm 9 \times 10^{-4}$	$0.8558 \pm 6 \times 10^{-4}$	$6 \times 10^{-5} \pm 4 \times 10^{-5}$	$6 \times 10^{-5} \pm 4 \times 10^{-5}$
gravity_gae	$0.9806 \pm 3 \times 10^{-4}$	$0.9753 \pm 4 \times 10^{-4}$	$0.069 \pm 6 \times 10^{-3}$	$0.101 \pm 5 \times 10^{-3}$
sourcetarget_gae	$0.99976 \pm 1 \times 10^{-5}$	$0.999736 \pm 8 \times 10^{-6}$	$0.077 \pm 4 \times 10^{-3}$	$0.147 \pm 7 \times 10^{-3}$
mlp_gae	$0.99315 \pm 1 \times 10^{-5}$	$0.99409 \pm 1 \times 10^{-5}$	$0.181 \pm 7 \times 10^{-3}$	$0.299 \pm 7 \times 10^{-3}$
digae	$0.9978 \pm 3 \times 10^{-4}$	$0.998 \pm 3 \times 10^{-4}$	$0.035 \pm 6 \times 10^{-3}$	$0.068 \pm 1 \times 10^{-2}$
magnet	$0.989 \pm 1 \times 10^{-4}$	$0.99076 \pm 3 \times 10^{-5}$	$0.151 \pm 1 \times 10^{-2}$	$0.26 \pm 2 \times 10^{-2}$

TABLE 1. Results for various graph autoencoder models.

11.1.5. *Discussion.* We observe consistently high *General* AUC and AUPRC scores, close to 1. These high values are expected, as similar neural architectures have demonstrated strong performance in graphs of comparable size [13]. The high ratio of existing to non-existing edges in the implication graph (approximately 37%) likely contributes to the near-perfect *General* AUC and AUPRC scores. For context, benchmark datasets such as Cora and Citeseer (e.g., directed and undirected) contain fewer than 1% of all possible edges.

Interestingly, the GAE model, though structurally unable to distinguish edge direction, performs well on the *General* task (if we consider AUC and AUPRC only). This experimentally confirms the need for including *Directional* and *Bidirectional* test sets, which allow comprehensive evaluation across all facets of Directed Link Prediction (DLP).

All other models demonstrate high AUC and AUPRC scores across the *General*, *Directional*, and *Bidirectional* test sets, indicating strong predictive capabilities even when accounting for directionality and bidirectionality.

Notably, the `mlp_gae` and `magnet` models also achieve competitive scores in MRR and Hits@K metrics.

11.1.6. *Conclusions.* We evaluated the performance of six different graph autoencoders on a Directed Link Prediction (DLP) task. By adopting a multi-task evaluation framework, we assessed the models comprehensively across various aspects of DLP. All models demonstrated strong performance on AUC and AUPRC metrics, with some also achieving high scores on MRR and Hits@K.

Node features were represented using one-hot encodings, meaning that the models received no explicit information about the equations represented by the nodes. Instead, they relied entirely on the topological structure encoded during training. This approach aligns with previous research suggesting that one-hot encodings can promote asymmetric embeddings [34]. However, future experiments could explore alternative representations, such as encoding the equations with text-based embeddings like Word2Vec, to potentially enhance the models' understanding of the nodes' semantic content.

In summary, our findings highlight the adaptability and robustness of graph autoencoders for DLP tasks in dense, directed graphs. This approach not only demonstrates robustness

in predicting directed links but also suggests promising potential for future improvements through enhanced feature representations, thereby advancing the capabilities of link prediction in complex mathematical knowledge graphs.

12. USER INTERFACE

Describe visualizations and explorer tools: Equation Explorer, Finite Magma Explorer, Graphiti, ...

13. STATISTICS AND EXPERIMENTS

TODO: Expand this sketch

Data analysis of the implication graph:

- Mention the long chain $2 \Rightarrow 5 \Rightarrow 2499 \Rightarrow 2415 \Rightarrow 238 \Rightarrow 2716 \Rightarrow 28 \Rightarrow 2973 \Rightarrow 270 \Rightarrow 3 \Rightarrow 3715 \Rightarrow 375 \Rightarrow 359 \Rightarrow 4065 \Rightarrow 1$ (discussed here).
- What are the "most difficult" implications?
- Is there a way to generate a standard test set of implication problems of various difficulty levels? Can one then use this to benchmark various automated and semi-automated methods? Challenge: how does one automatically assign a difficulty level to a given (anti-)implication?

See this for a preliminary data analysis of the impact of equation size and the number of variables.

Analyze the implication graph and discuss test sets of implication problems for benchmarking theorem provers. Challenge: How can one automatically assign a difficulty level to an implication?

14. DATA MANAGEMENT

TODO: expand this sketch

Describe how data was handled during the project and how it will be managed going forward.

15. REFLECTIONS

TODO: expand this sketch

Testimonies from individual participants (but perhaps this is more suited for a retrospective paper). Utilize the thoughts and reflections thread.

Automation often overtook the rate of human progress, for instance in developing metatheorems. Does this create an opportunity cost? Raised as a possible discussion point here by Zoltan Kocsis.

16. CONCLUSIONS AND FUTURE DIRECTIONS

TODO: Expand this sketch

Insights learned, future speculation. Utilize the discussion on future directions. Some ideas from that list:

- A database of triple implications (EquationX and EquationY imply EquationZ) - see also this discussion.
- Are there any equational laws that have no non-trivial finite models, but have surjective models?
- Can we produce interesting examples of irreducible implications EquationX \rightarrow EquationY (no intermediate EquationZ can interpose)
- Degree of satisfiability results, e.g., if a central groupoid obeys the natural central groupoid axiom 99% of the time, is it a natural central groupoid?
- Kisielewicz's question: does 5093 have any infinite models?
- "Insight mining" the large corpus of automated proofs that have been generated.
- How machine-learnable is the implication graph? (See AI/ML section)

ACKNOWLEDGMENTS

Thanks to Claudio Moroni for the exploration of directed link prediction on the implication graph using GNN autoencoders described in Section 11.1.

We are also grateful to the many additional participants to the Equational Theories Project for their numerous comments and encouragement, including but not restricted to Edward van de Meent, ...

APPENDIX A. NUMBERING SYSTEM

In this section we record the numbering conventions we use for equational laws.

For this formal definition we use the natural numbers $0, 1, 2, \dots$ to represent and order indeterminate variables; however, in the main text, we use the symbols $x, y, z, w, u, v, r, s, t$ instead (and do not consider any laws with more than eight variables).

We require to extend the ordering on indeterminate variables to a well-ordering on words w in the free magma generated by these variables by declaring $w > w'$ if one of the following holds:

- w has a larger order than w' .
- $w = w_1 \diamond w_2$ and $w' = w'_1 \diamond w'_2$ have the same order $n \geq 1$ with $w_1 > w'_1$.
- $w = w_1 \diamond w_2$ and $w' = w'_1 \diamond w'_2$ have the same order $n \geq 1$ with $w_1 = w'_1$ and $w_2 > w'_2$.

Thus for instance

$$0 < 1 < 0 \diamond 0 < 0 \diamond 1 < 1 \diamond 0$$

and

$$1 \diamond 1 < 0 \diamond (0 \diamond 0) < (0 \diamond 0) \diamond 0.$$

We similarly place a well-ordering on equational laws $w_1 \simeq w_2$ by declaring $w_1 \simeq w_2 > w'_1 \simeq w'_2$ if one of the following holds: as follows:

- $w_1 \simeq w_2$ has a longer order than $w'_1 \simeq w'_2$.
- If $w_1 \simeq w_2$ has the same order as $w'_1 \simeq w'_2$, and $w_1 > w'_1$.
- If $w_1 \simeq w_2$ has the same order as $w'_1 \simeq w'_2$, $w_1 = w'_1$, and $w_2 > w'_2$.

Two equational laws are equivalent if one can be obtained from another by some combination of relabeling the variables and applying the symmetric law $w_1 \simeq w_2 \iff w_2 \simeq w_1$. For instance, $(0 \diamond 1) \diamond 2 \simeq 1$ is equivalent to $0 \simeq (1 \diamond 0) \diamond 2$. We then replace every equational law with their minimal element in their equivalence class, which can be viewed as the *normal form* for that law; for instance, the normal form of $(0 \diamond 1) \diamond 2 \simeq 1$ would be $0 \simeq (1 \diamond 0) \diamond 2$. Finally, we eliminate any law of the form $w \simeq w$ other than $0 \simeq 0$. We then number the remaining equations $E1, E2, \dots$. For instance, $E1$ is the trivial law $0 \simeq 0$, $E2$ is the constant law $0 \simeq 1$, $E3$ is the idempotent law $0 \simeq 0 \diamond 0$, and so forth. Lists and code for generating these equations, or the equation number attached to a given equation, can be found at the ETP repository.

The number of equations in this list of order $n = 0, 1, 2, \dots$ is given by

$$2, 5, 39, 364, 4284, 57882, 888365, \dots$$

(<https://oeis.org/A376640>). The number can be computed to be

$$C_{n+1}B_{n+2}/2$$

if n is odd, 2 if $n = 0$, and

$$(C_{n+1}B_{n+2} + C_{n/2}(2D_{n+2} - B_{n+2}))/2 - C_{n/2}B_{n/2+1}$$

if $n > 2$ is even, where C_n, B_n are the Catalan and Bell numbers, and D_n is the number of partitions of $[n]$ up to reflection, which for $n = 0, 1, 2, \dots$ is

$$1, 1, 2, 4, 11, 32, 117, \dots$$

(<https://oeis.org/A103293>). A proof of this claim can be found in the ETP blueprint. In particular, there are 4694 equations of order at most 4.

Below we record some specific equations appearing in this paper, using the alphabet x, y, z, w, u, v in place of $0, 1, 2, 3, 4, 5, \dots$ for readability.

(E1)	$x \simeq x$	(Trivial law)
(E2)	$x \simeq y$	(Singleton law)
(E3)	$x \simeq x \diamond x$	(Idempotent law)
(E4)	$x \simeq x \diamond y$	(Left-absorptive law)
(E5)	$x \simeq y \diamond x$	(Right-absorptive law)
(E10)	$x \simeq x \diamond (y \diamond x)$	
(E23)	$x \simeq (x \diamond x) \diamond x$	
(E41)	$x \diamond x \simeq y \diamond z$	
(E43)	$x \diamond y \simeq y \diamond x$	(Commutative law)
(E46)	$x \diamond y \simeq z \diamond w$	(Constant law)
(E151)	$x \simeq (x \diamond x) \diamond (x \diamond x)$	
(E168)	$x \simeq (y \diamond x) \diamond (x \diamond z)$	(Central groupoid law)
(E206)	$x \simeq (x \diamond (x \diamond y)) \diamond y$	
(E327)	$x \diamond y \simeq x \diamond (y \diamond z)$	
(E378)	$x \simeq (x \diamond y) \diamond y$	
(E395)	$x \diamond y \simeq (z \diamond x) \diamond y$	
(E543)	$x \simeq y \diamond ((z \diamond (x \diamond (y \diamond z))))$	(Tarski's axiom)
(E854)	$x \simeq x \diamond ((y \diamond z) \diamond (x \diamond z))$	
(E1117)	$x \simeq y \diamond ((y \diamond (x \diamond z)) \diamond z)$	
(E1286)	$x \simeq y \diamond (((x \diamond y) \diamond x) \diamond y)$	
(E1485)	$x \simeq (y \diamond x) \diamond (x \diamond (z \diamond y))$	
(E1571)	$x \simeq (y \diamond z) \diamond (y \diamond (x \diamond z))$	(Exp. 2 abelian groups)
(E1648)	$x \simeq (x \diamond y) \diamond ((x \diamond y) \diamond y)$	
(E1689)	$x \simeq (y \diamond x) \diamond ((x \diamond z) \diamond z)$	
(E2441)	$x \simeq (x \diamond ((x \diamond x) \diamond x)) \diamond x$	
(E3316)	$x \diamond y \simeq x \diamond (y \diamond (x \diamond y))$	
(E4482)	$x \diamond (y \diamond y) = (y \diamond y) \diamond x$	
(E4512)	$(x \diamond y) \diamond z \simeq x \diamond (y \diamond z)$	(Associative law)
(E4531)	$x \diamond (y \diamond z) \simeq (y \diamond z) \diamond x$	
(E345169)	$x \simeq (y \diamond ((x \diamond y) \diamond y)) \diamond (x \diamond (z \diamond y))$	(Sheffer stroke)
(E42323216)	$x \simeq y \diamond (((y \diamond y) \diamond x) \diamond z)$	(Division in groups)
	$\diamond (((y \diamond y) \diamond y) \diamond z))$	

APPENDIX B. AUTHOR CONTRIBUTIONS

In this section we list the authors of this paper, grant support, and their contributions to this project, using the following standard CRediT categories:

- (1) Conceptualization: Ideas; formulation or evolution of overarching research goals and aims.
- (2) Data Curation: Management activities to annotate (produce metadata), scrub data and maintain research data (including software code, where it is necessary for interpreting the data itself) for initial use and later reuse.
- (3) Formal Analysis: Application of statistical, mathematical, computational, or other formal techniques to analyze or synthesize study data.
- (4) Funding Acquisition: Acquisition of the financial support for the project leading to this publication.
- (5) Investigation: Conducting a research and investigation process, specifically performing the experiments, or data/evidence collection.
- (6) Methodology: Development or design of methodology; creation of models.
- (7) Project Administration: Management and coordination responsibility for the research activity planning and execution.
- (8) Resources: Provision of study materials, reagents, materials, patients, laboratory samples, animals, instrumentation, computing resources, or other analysis tools.
- (9) Software: Programming, software development; designing computer programs; implementation of the computer code and supporting algorithms; testing of existing code components.
- (10) Supervision: Oversight and leadership responsibility for the research activity planning and execution, including mentorship external to the core team.
- (11) Validation: Verification, whether as a part of the activity or separate, of the overall replication/reproducibility of results/experiments and other research outputs.
- (12) Visualization: Preparation, creation and/or presentation of the published work, specifically visualization/data presentation.
- (13) Writing - Original Draft Preparation: Creation and/or presentation of the published work, specifically writing the initial draft (including substantive translation).
- (14) Writing – Review and Editing: Preparation, creation and/or presentation of the published work by those from the original research group, specifically critical review, commentary or revision – including pre- or post-publication stages.

- ...

- Pietro Monticone, Department of Mathematics, University of Trento, pietro.monticone@studenti.uni-trento.it: Data Curation, Formal Analysis, Project Administration, Resources, Software, Validation, Writing - Original Draft Preparation, Writing - Review and Editing.

- Terence Tao, Department of Mathematics, UCLA, tao@math.ucla.edu: Conceptualization, Data Annotation, Investigation, Methodology, Project Administration, Writing - Original Draft Preparation, Writing - Review and Editing. Supported by NSF grant DMS-2347850.

- ...

REFERENCES

- [1] A. K. Austin. A note on models of identities. *Proc. Amer. Math. Soc.*, 16:522–523, 1965.
- [2] Franz Baader and Tobias Nipkow. *Term rewriting and all that*. Cambridge University Press, Cambridge, 1998.
- [3] Alexander Bentkamp, Jasmin Blanchette, Visa Nummelin, Sophie Tournet, Petar Vukmirovic, and Uwe Waldmann. Mechanical mathematicians. *Commun. ACM*, 66(4):80–90, 2023.
- [4] George M Bergman. The diamond lemma for ring theory. *Advances in Mathematics*, 29(2):178–218, 1978.
- [5] Joshua Clune, Yicheng Qian, Alexander Bentkamp, and Jeremy Avigad. Duper: A proof-producing superposition theorem prover for dependent type theory. In Yves Bertot, Temur Kutsia, and Michael Norrish, editors, *15th International Conference on Interactive Theorem Proving, ITP 2024, September 9-14, 2024, Tbilisi, Georgia*, volume 309 of *LIPICs*, pages 10:1–10:20. Schloss Dagstuhl - Leibniz-Zentrum für Informatik, 2024.
- [6] Leonardo Mendonça de Moura and Nikolaj S. Bjørner. Efficient e-matching for SMT solvers. In Frank Pfenning, editor, *Automated Deduction - CADE-21, 21st International Conference on Automated Deduction, Bremen, Germany, July 17-20, 2007, Proceedings*, volume 4603 of *Lecture Notes in Computer Science*, pages 183–198. Springer, 2007.
- [7] Leonardo Mendonça de Moura and Nikolaj S. Bjørner. Z3: an efficient SMT solver. In C. R. Ramakrishnan and Jakob Rehof, editors, *Tools and Algorithms for the Construction and Analysis of Systems, 14th International Conference, TACAS 2008, Held as Part of the Joint European Conferences on Theory and Practice of Software, ETAPS 2008, Budapest, Hungary, March 29-April 6, 2008. Proceedings*, volume 4963 of *Lecture Notes in Computer Science*, pages 337–340. Springer, 2008.
- [8] Gabriel Ebner, Sebastian Ullrich, Jared Roesch, Jeremy Avigad, and Leonardo de Moura. A metaprogramming framework for formal verification. *Proc. ACM Program. Lang.*, 1(ICFP):34:1–34:29, 2017.
- [9] Trevor Evans. Products of points—some simple algebras and their identities. *Amer. Math. Monthly*, 74:362–372, 1967.
- [10] Yifan Feng, Haoxuan You, Zizhao Zhang, Rongrong Ji, and Yue Gao. Hypergraph neural networks. *Proceedings of the AAAI Conference on Artificial Intelligence*, 33(01):3558–3565, July 2019.
- [11] Timothy Gowers and Michael Nielsen. Massively collaborative mathematics. *Nature*, 461(7266):879–881, October 2009.
- [12] Graham Higman and B. H. Neumann. Groups as groupoids with one law. *Publ. Math. Debrecen*, 2:215–221, 1952.
- [13] Thomas N. Kipf and Max Welling. Variational graph auto-encoders. 2016.
- [14] A. Kisielewicz. Austin identities. *Algebra Universalis*, 38(3):324–328, 1997.
- [15] M. Kivela, A. Arenas, M. Barthélemy, J. P. Gleeson, Y. Moreno, and M. A. Porter. Multilayer networks. *Journal of Complex Networks*, 2(3):203–271, July 2014.
- [16] Donald E. Knuth. Notes on central groupoids. *J. Combinatorial Theory*, 8:376–390, 1970.
- [17] Donald E. Knuth and Peter B. Bendix. Simple word problems in universal algebras. In *Computational Problems in Abstract Algebra (Proc. Conf., Oxford, 1967)*, pages 263–297. Pergamon, Oxford-New York-Toronto, Ont., 1970.
- [18] Thomas Koehler, Andrés Goens, Siddharth Bhat, Tobias Grosser, Phil Trinder, and Michel Steuwer. Guided equality saturation. *Proc. ACM Program. Lang.*, 8(POPL):1727–1758, 2024.
- [19] Georgios Kollias, Vasileios Kalantzis, Tsuyoshi Idé, Aurélie Lozano, and Naoki Abe. Directed graph auto-encoders. 2022.
- [20] Laura Kovács and Andrei Voronkov. First-order theorem proving and vampire. In Natasha Sharygina and Helmut Veith, editors, *Computer Aided Verification - 25th International Conference, CAV 2013, Saint Petersburg, Russia, July 13-19, 2013. Proceedings*, volume 8044 of *Lecture Notes in Computer Science*, pages 1–35. Springer, 2013.
- [21] Daniel Kroening and Ofer Strichman. *Decision Procedures - An Algorithmic Point of View, Second Edition*. Texts in Theoretical Computer Science. An EATCS Series. Springer, 2016.
- [22] André Kündgen, Gregor Leander, and Carsten Thomassen. Switchings, extensions, and reductions in central digraphs. *J. Combin. Theory Ser. A*, 118(7):2025–2034, 2011.

- [23] Juanhui Li, Harry Shomer, Haitao Mao, Shenglai Zeng, Yao Ma, Neil Shah, Jiliang Tang, and Dawei Yin. Evaluating graph neural networks for link prediction: Current pitfalls and new benchmarking. 2023.
- [24] Jannis Limperg and Asta Halkjær From. Aesop: White-box best-first proof search for lean. In Robbert Krebbers, Dmitriy Traytel, Brigitte Pientka, and Steve Zdancewic, editors, *Proceedings of the 12th ACM SIGPLAN International Conference on Certified Programs and Proofs, CPP 2023, Boston, MA, USA, January 16-17, 2023*, pages 253–266. ACM, 2023.
- [25] W. McCune. Prover9 and mace4. <http://www.cs.unm.edu/~mccune/prover9/>, 2005–2010.
- [26] William McCune. Solution of the Robbins problem. *J. Automat. Reason.*, 19(3):263–276, 1997.
- [27] William McCune. Single axioms: with and without computers. In *Computer mathematics (Chiang Mai, 2000)*, volume 8 of *Lecture Notes Ser. Comput.*, pages 83–89. World Sci. Publ., River Edge, NJ, 2000.
- [28] William McCune, Robert Veroff, Branden Fitelson, Kenneth Harris, Andrew Feist, and Larry Wos. Short single axioms for Boolean algebra. *J. Automat. Reason.*, 29(1):1–16, 2002.
- [29] Ralph McKenzie. On spectra, and the negative solution of the decision problem for identities having a finite nontrivial model. *J. Symbolic Logic*, 40:186–196, 1975.
- [30] N. S. Mendelsohn and R. Padmanabhan. Minimal identities for Boolean groups. *J. Algebra*, 34:451–457, 1975.
- [31] C. A. Meredith and A. N. Prior. Equational logic. *Notre Dame J. Formal Logic*, 9:212–226, 1968.
- [32] J. D. Phillips and Petr Vojtěchovský. The varieties of loops of Bol-Moufang type. *Algebra Universalis*, 54(3):259–271, 2005.
- [33] Marcus Rossel and Andrés Goens. Bridging syntax and semantics of lean expressions in e-graphs. *arXiv preprint arXiv:2405.10188*, 2024.
- [34] Guillaume Salha, Stratis Limnios, Romain Hennequin, Viet Anh Tran, and Michalis Vazirgiannis. Gravity-inspired graph autoencoders for directed link prediction. 2019.
- [35] Alfred Tarski. Ein beitrag zur axiomatik der abelschen gruppen. *Fundamenta Mathematicae*, 30(1):253–256, 1938.
- [36] Max Willsey, Chandrakana Nandi, Yisu Remy Wang, Oliver Flatt, Zachary Tatlock, and Pavel Panchekha. egg: Fast and extensible equality saturation. *Proc. ACM Program. Lang.*, 5(POPL):1–29, 2021.
- [37] Stephen Wolfram. The physicalization of metamathematics and its implications for the foundations of mathematics, Mar 2022.
- [38] Michihiro Yasunaga, Antoine Bosselut, Hongyu Ren, Xikun Zhang, Christopher D Manning, Percy Liang, and Jure Leskovec. Deep bidirectional language-knowledge graph pretraining. 2022.
- [39] Xikun Zhang, Antoine Bosselut, Michihiro Yasunaga, Hongyu Ren, Percy Liang, Christopher D. Manning, and Jure Leskovec. Greaselm: Graph reasoning enhanced language models for question answering. 2022.
- [40] Xitong Zhang, Yixuan He, Nathan Brugnone, Michael Perlmutter, and Matthew Hirn. Magnet: A neural network for directed graphs. 2021.